# Transformers

Transformers in natural language processing (NLP) are a type of deep learning module that uses self-attention mechanisms to analyze and process natural language data. They are encoder decoder models that can be used for many applications, including machine translation

## Attention Mechanism

Parallely we cannot send all the words in a sentence

## Transformers

Using self attention to parallely send all the words

---> Architecture